

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

Sleep Queue Management

Inventor:
Bor-Ming Hsieh

ATTORNEY'S DOCKET NO. MS1-742US

RELATED APPLICATIONS

This patent application is related to U.S. patent application serial no. _____, filed on _____, and titled "Run Queue Management", which is hereby incorporated by reference.

TECHNICAL FIELD

The following description relates to operating systems. More particularly, the following description relates to sleep queue management.

BACKGROUND

Real-time performance is essential for time-critical responses required in high-performance embedded applications such as telecommunications switching equipment, medical monitoring equipment, space navigation and guidance applications, and the like. Such applications must deliver responses within specified time parameters in real-time.

Real-time performance is provided by operating systems that use thread-scheduling mechanisms to schedule threads for execution on a thread priority basis. (A thread is basically a path of execution through a computer program application). For example, the Microsoft WINDOWS CE ® operating system provides two-hundred-fifty-six (256) thread priority levels, wherein threads of higher priority are executed before threads of lower priority are executed. Threads of equal priority are executed in a first-in-first-out round-robin fashion. For example, thread A runs, then thread B runs, followed by thread C, and back to thread A.

1 Thread scheduling mechanisms typically store threads in a "run queue" for
2 subsequent execution. Fig. 1 illustrates a traditional run queue 100 that is
3 implemented as a linked list. The threads 102 in the run queue are sorted based on
4 respective thread priorities. For example, Threads 102-A1 through 102-AJ have
5 respective thread priorities of zero (0), threads 102-B1 through 102-BK have
6 respective thread priorities of one (1), and the like. There can be any number of
7 threads 102 in the run queue.

8 A thread that is currently executing may be preempted by another thread, or
9 the thread itself may "yield" its access to the processor. Both of these terms refer
10 to the thread being placed into a "non-executing" state. For example, an operating
11 system may put a thread to sleep, or preempt the thread to allow a different thread
12 with a higher priority to execute. In another example, the thread itself may "yield"
13 its access to the processor to wait for the occurrence of a particular event such as
14 the elapse of a predetermined amount of time, or the like, before continuing
15 execution.

16 Regardless of whether a thread is preempted by another program or whether
17 the thread itself yields its access to the processor, the system's thread scheduling
18 mechanism typically stores the preempted, or yielding thread into a run queue or
19 sleep queue. (More particularly, a reference to the thread is generally stored in the
20 run queue or sleep queue). (Although the thread scheduling mechanism may or
21 may not be part of the operating system, the terms thread scheduling mechanism,
22 operating system, and the like, are often used interchangeably in this description to
23 describe a system's thread scheduling aspects). When a thread's specified sleep
24 time has expired, the scheduling mechanism "wakes-up" the thread by removing
25

the thread from the sleep queue and inserting the thread into the run queue for subsequent execution.

Fig. 2 illustrates a traditional single-dimension sleep queue 200 that is implemented as a linked list. (A traditional sleep queue may also be implemented as a "heap" data structure). For purposes of this description, a sleep queue 200 is any queue for storing any number of threads that are sorted based on time. In this example, the threads 202 in the sleep queue are sorted in a single dimension based on thread wake-up time and thread priority within a particular wake-up time. For example, thread 202-1 has a wake-up time of five (5) milliseconds (ms) and threads 202-2 and 202-3 have respective wake-up times of ten (10) milliseconds. Threads that have the same sleep time are sorted based on priority in a round robin fashion. For example, thread 202-2 has a wake-up time of 10 ms with a thread priority of 0 (in this example, the highest thread priority), and thread 202-... has a wake-up time of 10 ms with a thread priority of 5 (a lower priority than a thread priority of 0). In this manner, the threads in the traditional sleep queue are sorted with respect to one-another in a single dimension.

As discussed above, a thread may be preempted by another thread for any number of reasons. One significant reason that a thread may be preempted is so that the operating system, or thread scheduling mechanism can determine if there are any threads of higher priority that need to be executed. Part of this determination, and another significant reason in and of itself, is the operating system may scan the threads stored/referenced in the sleep queue to determine if any need to be woken-up for execution (e.g., inserted into the run queue). Real-time operating systems typically preempt all other threads from executing at predetermined periodic time intervals to perform such thread management.

1 Thread scheduling mechanisms typically use a hardware timer to produce a
2 system tick to determine a maximum amount of time, or "quantum" that a thread
3 can execute in the system without being preempted. A system tick is a rate at
4 which a hardware timer interrupt is generated and serviced by an operating
5 system. When the timer fires, the thread scheduling mechanism will schedule a
6 new thread for execution if one is ready.

7 Significantly, an operating system requires exclusive access to a processor
8 during certain thread scheduling procedures such as during sleep queue thread
9 removal procedures and during run queue thread insertion procedures. The
10 operating system uses its system-exclusive access: (a) to remove threads from the
11 sleep queue at or as close as possible to each respective thread's specified wake-up
12 time for subsequent insertion into the run queue; and, (b) to insert each thread
13 removed from the sleep queue into the run queue for execution.

14 The number of threads to be woken-up at any one time could be any
15 number of threads such as one thread, two threads, or one hundred threads. The
16 more threads that need to be removed from the sleep queue for insertion into the
17 run queue, the greater the amount time is that an operating system requires system-
18 exclusive access to the processor. This system-exclusive access is directly
19 controlled by the operating system and cannot typically be preempted by any other
20 thread.

21 The non-deterministic and non-preemptable nature of traditional sleep
22 queue thread removal and run queue thread insertion procedures creates a number
23 of significant problems. One problem, for example, is that an operating system
24 cannot typically be guaranteed to schedule other threads within predetermined
25 time parameters because of such non-deterministic thread management techniques.

1 This means that a preempted thread (a thread that was executed but that was
2 blocked during sleep queue thread removal) won't execute again for an unknown
3 amount of time. The respective wake-up times of one or all of the threads that that
4 need to be removed from a sleep queue at any one moment in time may have
5 already long passed before they are removed and inserted into the run queue.
6 Analogously, by the time a thread that is inserted into the run queue gets executed,
7 the thread's purpose or the event that the thread is responding to may have passed
8 long ago.

9 Accordingly, traditional sleep queue thread removal and run queue thread
10 insertion procedures do not typically allow an operating system to schedule other
11 threads for execution within deterministic/predetermined time parameters.

12 13 SUMMARY

14 Various implementations of the described subject matter provide for the
15 management of a multi-dimensional sleep queue, such that a group of threads with
16 a same wake-up time are removed from the multi-dimensional sleep queue in a
17 deterministic amount of time independent of the number of threads in the group.
18 Moreover, new threads are inserted into the multi-dimensional sleep queue in a
19 manner that allows other processes to execute during the thread insertion process.
20 Thus, whether inserting threads into a sleep queue, or whether removing threads
21 from the sleep queue, the described subject matter allows an operating system to
22 schedule other threads for execution within deterministic/predetermined time
23 parameters.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a block diagram that illustrates aspects of a traditional run queue used by an operating system to schedule threads for execution.

Fig. 2 is a block diagram that illustrates aspects of a traditional thread sleep queue used by an operating system to allow threads to sleep for a specified amount of time before execution is resumed.

Fig. 3 is a block diagram that shows an exemplary multi-dimensional sleep queue that is used by a thread scheduling mechanism to allow an operating system to schedule other threads for execution within deterministic/predetermined time parameters.

Fig. 4 is a flowchart that illustrates an exemplary procedure to insert a thread into a multi-dimensional sleep queue, such that multiple threads that have a same thread wake-up time can be removed from the sleep queue in a deterministic amount of time.

Fig. 5 is a flowchart that shows an exemplary procedure to determine/establish a thread insertion point in a multi-dimensional sleep queue in a manner that allows an operating system to schedule other threads for execution within deterministic time parameters.

Fig. 6 is a flowchart illustrating further aspects of an exemplary optimized procedure to insert a thread into a multi-dimensional sleep queue. Specifically, the optimized procedure uses a multi-dimensional atomic walk procedure to identify a position in a multi-dimensional sleep queue to insert the thread.

Fig. 7 is a flowchart that shows further aspects of an exemplary multi-dimensional atomic walk procedure to identify a position in a multi-dimensional sleep queue to insert a thread. Specifically, Fig. 7 illustrates use of a last examined

thread to identify a start position in a multi-dimensional sleep queue to begin a search for a new thread insertion point.

Fig. 8 is a flowchart that shows further aspects of a multi-dimensional atomic walk procedure to insert a new thread into a multi-dimensional sleep queue. In particular Fig. 8 shows how a last examined node may be used to identify an insertion point in the sleep queue based on the last examined node's wake-up time and priority.

Fig. 9 is a flowchart that illustrates further aspects of an exemplary optimized procedure to insert a new thread into a multi-dimensional sleep queue. In particular, Fig. 9 shows how a last examined node is used to identify an insertion point in the sleep queue when the last examined thread is the last thread in one of the multiple dimensions.

Fig. 10 is a flowchart that shows an exemplary procedure to remove a group of threads from a sleep queue in a deterministic amount of time that is independent of the number of threads that need to be removed.

Fig. 11 is a block diagram that shows aspects of an exemplary run queue that allows a plurality of threads to be associated with the run queue in a deterministic amount of time.

Fig. 12 is a flowchart that shows aspects of an exemplary procedure to remove a thread from run queue.

Fig. 13 is a flowchart diagram that illustrates aspects of an exemplary procedure to insert a plurality of threads into a run queue in a determinative amount of time that is independent of the number of threads associated with the run queue at any one time.

Fig. 14 is block diagram that shows aspects of a suitable computing environment wherein an exemplary system and procedure to manage a run queue may be implemented.

DETAILED DESCRIPTION

The following description sets forth various implementations of subject matter to manage a sleep queue that incorporates elements recited in the appended claims. The implementations are described with specificity in order to meet statutory requirements. However, the description itself is not intended to limit the scope of this patent. Rather, the inventor has contemplated that the claimed subject matter might also be embodied in other ways, to include different elements or combinations of elements similar to the ones described in this document, in conjunction with other present or future technologies.

Exemplary Multi-Dimensional Sleep Queue

Fig. 3 is a block diagram that shows an exemplary two-dimensional (2-D) sleep queue 300. Sleep queue 300 is used to keep track of threads that are no longer executing for a specified amount of time, or “put to sleep”. When a thread is interrupted and put to sleep, a reference to that thread is inserted in the sleep queue. Each thread has parameters that are used to determine when it should be taken out of the sleep queue and therefore “resumed”. In this implementation, these parameters include a wake-up time and a thread priority.

For example, within the sleep queue 300, each thread is represented as a node or data field 304. Thus, the sleep queue comprises a plurality of such nodes or data fields. Each node represents a respective thread of execution in a computer

program application. Each node comprises a respective thread wake-up time and a respective thread priority ("PRI"). For example, node 304-1 includes a thread WT of five (5) milliseconds (ms), and a thread PRI of ten (10). (In this implementation, a thread PRI of 0 is a higher thread PRI than one (1), which in turn is a higher thread PRI than two (2), and the like).

For purposes of removing a group of threads from the sleep queue 300 in a deterministic amount of time, the nodes 304 are sorted in two dimensions: first with respect to their wake-up times and then with respect to their thread priorities. Fig. 3 shows the results of such sorting, in which nodes are sorted left-to-right in order of increasing wake-up time values and top-to-bottom in order of increasing PRI values. This produces a two-dimensional array having a number of columns equal to the number of unique wake-up time values possessed by the nodes contained in the sleep queue. Each column contains all nodes having a particular wake-up time value. The columns are sorted in order of increasing wake-up time values. Within each column, the nodes are sorted in order of their PRI values.

For example, a first-dimension row or data field 302 includes nodes 304-1 through 304-5, having the lowest-priority node of each unique wake-up time value, sorted in order of increasing wake-up time value. Node 304-1 has a thread wake-up time of 5, node 304-2 has a thread wake-up time of 7, node 304-3 has a thread wake-up time of 20, and the like. Each node is ordered in the first dimension such that a thread that is ready to be "woken-up" is removed by a sleep queue management procedure before the procedure removes a thread whose "wake-up" time has not arrived. There can be any number of nodes in the first dimension data field. Each node 304 in the first dimension 302 has a different thread wake-up time as compared to other nodes in the first dimension. First

1 dimension nodes are not ordered, or sorted with respect to one-another based on
2 thread priority.

3 The sorting also results in a plurality of second-dimension columns or data
4 fields 308, wherein nodes 304 are ordered with respect to one-another based on the
5 their respective thread PRIs. For example, a second subset of the nodes is
6 represented by second dimension data field 308, which in this example comprises
7 nodes 304-2 through 304-8. There can be any number of nodes in the second
8 dimension data field. Moreover, there can be any number of second dimension
9 data fields. To illustrate this, consider that there is a second dimension data
10 field 308-2, which comprises nodes 304-4 through 304-10.

11 Each node 304 represents not only a node in the first dimension 302, but
12 also represents a respective node in a second dimension 308. In other words, each
13 second dimension data field includes a single first dimension node. For example,
14 second dimension 308-1 includes node 304-1, which is in both the first dimension
15 and the second dimension. In yet another example, second dimension 308-2
16 includes node 304-4, which is sorted with respect to both the first dimension and
17 the second dimension.

18 As illustrated by directional arrow 310, each node 304 in a second
19 dimension 308 is sorted with respect to each other node in the second dimension
20 based on thread PRI within a particular thread wake-up time. For example, node
21 304-2 has a thread PRI of two (2), which is higher than node's 304-6 respective
22 thread priority of three (3), and the like. In contrast to the nodes 304 in the first
23 dimension 302, wherein each node has a different thread wake-up time, each
24 second dimension node has a same respective thread wake-up time. For example,
25 each node 304-2 through 304-8 has the same thread wake-up time of seven (7) ms.

Because each node 304 of the sleep queue 300 is sorted with respect to first and second dimensions, each node with a highest/higher priority in the first dimension 302 as compared to respective priorities of other nodes with a similar wake-up time is considered to be in a "primary position" with respect to the first and second 308 dimensions of the sleep queue. For example, referring to Fig. 3, nodes 304-1 through 304-5 can be considered to be in a "primary" position with respect to the first and second dimensions of the sleep queue.

Furthermore, a node with a lower priority as compared to one or more other nodes with a similar wake-up time, is considered to be in a secondary position with respect to the first and second dimensions of the sleep queue. For example, node 304-6 has a priority of 3, which in this example is lower than node's 304-2 priority of 2. (Note that each node has a similar wake-up time of 7 ms). Thus, node 304-2 is considered to be in a secondary position with respect to the first and second dimensions of the sleep queue.

To indicate sorted relationships between nodes, the nodes are advantageously maintained in linked lists. Specifically, nodes 304 may include one or more references to other nodes 304. For example, a first thread 304-1 includes a reference (not shown) such as a pointer reference to a next thread 304-2, which in turn comprises a reference to a subsequent thread 304-3. In yet another example, node 304-2 includes at least two such node references, a first reference to node 304-6 and a second reference to 304-3.

There are many ways of managing such node 304 references. For example, if a node does not currently reference a different node, a pointer reference in the node may be set to equal a null value. If during sleep queue management procedures the multi-dimensional sleep queue configuration changes to include

more or less threads/nodes, a node's corresponding pointer reference can be set accordingly to accurately reflect such changes.

This multi-dimensional arrangement of the sleep queue 300 makes it easy to determine the next thread to be resumed or group of threads to be removed from the sleep queue. In this example, the next thread to be removed is the thread with the earliest wake-up time and highest thread priority—node 304-1 at the top left. Furthermore, this arrangement has several other advantages relating to removing groups of threads from the sleep queue in a deterministic amount of time. For example, when the sleep queue 300 is implemented as a linked list, a single node (e.g. node 304) detach operation is used to remove one or more nodes in the first 302 and second dimensions 308 on a system timer tick.

To illustrate this, consider that a detach of node 304-2 from the multi-dimensional sleep queue results in the removal of a group of nodes (i.e., node 304-2, and nodes 304-6 through 304-8) in second dimension data field 308-1. Thus, each of the nodes in a second dimension data field, regardless of the number of nodes in the second dimension data field, is removed from the sleep queue at the same time.

Accordingly, the amount of time it takes to remove any number of nodes with a particular wake-up time is based on the amount of time it takes to detach a single node from the sleep queue. Thus, the multi-dimensional sleep queue provides for removing a group of nodes from the sleep queue in a bounded, or determinative amount of time. As described below, this group of nodes can be associated with the run queue in the amount of time that it takes to insert a single node into the run queue. Thus, use of sleep queue 300 during thread management

allows an operating system to schedule other threads for execution within deterministic/predetermined time parameters.

Exemplary Procedure to Manage a Multi-Dimensional Sleep Queue

Conventional techniques to insert threads into a single dimension sleep queue of Fig. 2 typically allow other processes/threads to execute during the sleep queue thread insertion process. This allows an operating system to schedule other threads for execution within predetermined time parameters during the sleep queue thread insertion process. (This is the case even though the entire amount of time that it takes to insert a thread into the sleep queue may be non-deterministic, or unbounded (undeterminable)). However, such conventional single dimension sleep queue thread insertion techniques are not designed to insert a thread into the inventive multi-dimensional sleep queue 300 of Fig. 3.

Fig. 4 is a flowchart that shows an exemplary procedure 400 to insert a thread into a multi-dimensional sleep queue 300 of Fig. 3. The procedure inserts threads into the multi-dimensional sleep queue such that multiple threads having a same thread wake-up time can be removed from the sleep queue in a determinative amount of time, and also such that the wake-time and priority sorted semantics of the multi-dimensional sleep queue are preserved.

At block 410, the procedure receives a new thread of execution to be inserted into a multi-dimensional sleep queue for a predetermined, thread specified amount of time. (See, the sleep queue 300 of Fig. 3). At some point in time, no threads may be stored in the sleep queue. Thus, at block 412, the procedure determines if the new thread will be the first thread in the sleep queue. If not, the procedure continues at block 416, which is described in greater detail below. If so,

1 at 414, the procedure inserts the new thread into the sleep queue as a first thread,
2 or "head node" of the sleep queue.

3 At block 416, the procedure 400 having already determined that the new
4 thread is not the first thread in the sleep queue (block 412), the procedure
5 establishes a thread insertion point in the multi-dimensional sleep queue for the
6 new thread. The thread insertion point is based on the new thread's specific wake-
7 up time and priority, as compared to the respective wake-up times and priorities of
8 threads that are already in the sleep queue. (Exemplary aspects to establish a
9 proper thread insertion point are described in greater detail below in reference to
10 Figs. 5 through 9). At block 418, the procedure introduces the new thread into the
11 sleep queue at the established insertion point.

12 Fig. 5 is a flowchart that illustrates an exemplary block 416 of the
13 procedure of Fig. 4 to establish an insertion point in a multi-dimensional sleep
14 queue. At block 510, the procedure determines if the new thread's specified wake-
15 up time is different as compared to each respective wake-up time of each of the
16 other threads in a first dimension of the multi-dimensional sleep queue. At
17 block 512, it having been determined that the new thread's wake-up time is unique
18 in the first dimension (block 510), the procedure sets the insertion point for the
19 new thread such that it is sorted into the first dimension (see data field 302 of Fig.
20 3) based on its unique wake-up time.

21 At block 514, having already determined that a different thread in the first
22 dimension has a same wake-up time as the new thread wake-up time (block 510),
23 the procedure 400 determines whether the new thread's priority is the same or
24 higher than a thread priority that corresponds to the different thread with the same
25 wake-up time.

1 At block 516, it having been determined that the new thread has a same or
2 higher priority as compared to the different thread (block 514), the procedure
3 establishes the insertion point within a second dimension of threads having a
4 similar wake-up time. If thread priorities are the same, the new thread insert point
5 is immediately before or after the different thread. In this implementation, if the
6 new thread's priority is higher or the same as the different thread's priority, the
7 insert point is immediately before the different thread.

8 At block 518, having already determined that a different thread in the first
9 dimension has a same wake-up time as the new thread wake-up time (block 510),
10 and having already determined that the new thread does not have a same or higher
11 priority than the different thread (block 514), the procedure 400 establishes the
12 insert point based on the new thread's lower priority within a second dimension of
13 nodes with the same wake-up time.

14 Fig. 6 is a flowchart that illustrates an exemplary optimized procedure to
15 insert a new thread into a multi-dimensional sleep queue 300 of Fig. 3. More
16 particularly, this implementation uses a multi-dimensional atomic walk to locate a
17 proper insertion point in the multi-dimensional sleep queue for a new thread. The
18 multi-dimensional atomic walk either starts searching for an insertion point at a
19 first thread (if there is one), or at a last examined node in the sleep queue.

20 If a status of a last examined node 302 has changed, the search will begin at
21 the start of the sleep queue. Such a change of node status comprises a
22 determination of whether the last examined node was already removed from the
23 sleep queue since it was last examined, or whether the last examined node was
24 moved from a primary position with respect to the first and second dimensions of
25 the sleep queue to a secondary position. (Primary and secondary positions with

1 respect to the first and second dimensions of the sleep queue are discussed in
2 greater detail above in reference to Fig. 3). As long the status of the last examined
3 thread has not changed, it is valid to begin the examination of threads in the sleep
4 queue with the last examined thread.

5 The search for the thread insertion point is performed on a thread-by-
6 thread, or "node-by-node" basis. A node-by-node basis means that an operating
7 system maintains system-exclusive access to a processor only for that amount of
8 time that is required to examine a single node to determine if the examined node
9 identifies a new thread's appropriate insertion point in a multi-dimensional sleep
10 queue. After a single node is examined, the operating system releases the system-
11 exclusive access to the processor. If yet another node needs to be examined to
12 identify the insert point for the new thread, then the operating system again grabs
13 system-exclusive access to the processor to examine a next thread/node (if any).
14 In this manner, in-between single node examinations, the operating system allows
15 the processor to execute threads that were pre-empted to allow the operating
16 system to perform the sleep queue scheduling mechanism. An optimized multi-
17 dimensional atomic walk thread insertion procedure 600 is now described.

18 At block 610 the procedure 600 determines if a new thread that is to be
19 inserted into the multi-dimensional sleep queue is a first thread. If so, at block 612
20 the procedure inserts the first thread into the sleep queue.

21 If the new thread is determined not to be the first thread (block 610), at
22 block 618, the procedure sets a last examined node/thread to reference, the
23 inserted first node (block 612). At block 620, the procedure preempts all other
24 threads from executing by grabbing system-exclusive access to the processor.
25

1 At block 622, the procedure 600 determines if a state of the last node has
2 changed. As discussed above, the last node's state changes if it has already been
3 removed from the sleep queue (e.g., already inserted into the run queue for
4 execution), or if the last node was moved from a primary position with respect to
5 the first and second dimensions of the sleep queue to a secondary position.

6 If the state of the last node has not changed (block 622), the procedure 600
7 continues at block 710 as shown in Fig. 7, which is described in greater detail
8 below. However, if the state has changed, at block 624, the operating system
9 releases the system-exclusive access to the processor (see, block 620). The
10 procedure continues at block 614, wherein the procedure determines if it is time
11 for a thread to be woken-up from the sleep queue. If so, the procedure ends.
12 Otherwise, the procedure continues as described above in reference to block 618.

13 Fig. 7 is a flowchart that shows further aspects of an exemplary
14 procedure 600 to insert a new thread into a multi-dimensional sleep queue to. At
15 block 710, the procedure determines if the last examined node/thread indicates an
16 insertion point in the sleep queue. (An exemplary methodology of block 710 to
17 determine if the last node identifies an insertion point is described in greater detail
18 below in reference to Fig. 8). If so, at block 712, the procedure inserts the new
19 node/thread into the multi-dimensional sleep queue at the indicated insert point.
20 At block 714, the procedure releases system-exclusive access to the processor
21 (see, block 620 of Fig. 6).

22 At block 716, it having been determined that the last examined thread/node
23 does not indicate an insertion point for the new thread in the sleep queue
24 (block 710), the procedure 600 sets the last examined node/thread to indicate a
25 next node in the sleep queue. At block 718, the procedure releases the system-

1 exclusive access to the processor (see, block 620 of Fig. 6). At block 720, the
2 procedure determines if a thread needs to be woken-up from the sleep queue. If
3 so, the procedure ends. Otherwise, the procedure continues at block 620 of Fig. 6,
4 which is described in greater detail above.

5 Fig. 8 is a flowchart that shows further aspect of a multi-dimensional
6 atomic walk procedure to insert a new thread into a multi-dimensional sleep
7 queue. In particular it shows how a last examined node may identify an insertion
8 point in the sleep queue. (See, block 710 of Fig. 7). At block 810, the procedure
9 determines if the new node/thread has an earlier wake-up time as compared to the
10 last examined node/thread. If so, the insert point is established based on the new
11 thread's earlier WT, such that the new thread will be removed from the sleep queue
12 before the last examined thread.

13 At block 814, the procedure 600 determines if the new node/thread has a
14 same wake-up time as compared to the last examined node/thread. If not, the
15 procedure continues at block 910 of Fig. 9, which is described in greater detail
16 below. However, if the new node/thread has a same wake-up time as compared to
17 the last examined node/thread, at block 816, the procedure determines if the new
18 node/thread has a same or higher priority as compared to the priority of the last
19 examined node/thread. If not, the procedure continues at block 910 of Fig. 9,
20 which is described in greater detail below.

21 If the new node/thread has a same wake-up time and a same or higher
22 priority as compared to the last examined node/thread (block 816), at block 818,
23 the procedure establishes the insert point based on the new thread's same or higher
24 priority as compared to the priority of the last examined node/thread, and based on
25 the similar wake-up time as the last examined node/thread. (Exemplary

methodology to perform block 818 is described in greater detail above in reference to block 516 of Fig. 5).

Fig. 9 is a flowchart that shows further aspect of a multi-dimensional atomic walk procedure to insert a new thread into a multi-dimensional sleep queue. In particular it shows how a last examined node that is a last node in a dimension may identify an insertion point in the sleep queue for a new thread. (See, block 710 of Fig. 7, and block 816 of Fig. 8). At block 910, the procedure determines if the last examined node indicates a next node (or is it null) in a particular dimensional direction of interest to indicate whether there are any other next nodes in the primary or secondary dimensions of the multi-dimensional sleep queue.

If the last examined node has the same wake-up time as the new node (this was already determined at block 814 of Fig. 8), then it has also already been determined that the new node has a lower priority than the last node (see, block 816 of Fig. 8). Thus, the procedure determines at block 910 whether the last examined node indicates a next node with respect to a node with a lower priority in the second dimension (e.g., a next node with a secondary position with respect to the first and second dimensions).

If the last examined node does not have the same wake-up time as the new node (this was already determined at block 814 of Fig. 8), then it has also already been determined that the new thread/node has a later wake-up time as compared to the last node (see, block 810 of Fig. 8). Thus, the procedure determines at block 910 whether the last examined node indicates a next node with respect to a node with a later wake-up time than the last node in the first dimension (e.g., a next node with a primary position with respect to the first and second dimensions).

At block 912, it having been determined that the last node is the last node in a dimensional direction of interest, establishes the insert point for the new thread such that the new node is the last node in that dimensional direction of interest and such that it will be removed from the multi-dimensional sleep queue after the last examined node.

At block 914, the last node not being the last node in a dimensional direction of interest, the procedure indicates that the new thread's insert point in the multi-dimensional sleep queue is not yet determinable. The procedure continues at block 710 of Fig. 7 as described in greater detail above.

Time-Deterministic Group Thread Removal from a Sleep Queue

Fig. 10 is a flowchart that shows an exemplary procedure 1000 to remove a group of threads from a sleep queue in a bounded, or deterministic amount of time. At block 1010, the procedure determines whether one or more respective thread specified wake-up times have expired. (A thread with an expired wake-up time must be removed from the sleep queue for subsequent insertion into a run queue for execution). At block 1012, having determined that one or more threads need to be removed from the sleep queue for insertion into the run queue (block 1010), the procedure removes the one or more threads from the sleep queue in the amount of time the processor takes to perform a single node detach operation. Thus, the group removal of threads from the multi-dimensional sleep queue is time-deterministic in nature.

Accordingly, the methodology shown in Figs. 5-10 provide for the removal of multiple threads from a sleep queue in a deterministic amount of time. Moreover, as will be discussed in greater detail below in reference to Figs. 11 and

12, the amount of time that it will take to associate the detached group of nodes with the run queue is also time-deterministic because it will only be that amount of time that it takes to associate a single node with the run queue. This is significant because any non-deterministic delay incurred by the operating system in providing program responses threatens the real-time aspects of a real-time operating system.

Exemplary Run Queue

Fig. 11 is a block diagram that shows aspects of an exemplary run queue 1100 for associating a group of nodes with the run queue in a deterministic amount of time. The run queue comprises a first data field 1102 includes a first plurality of threads 1104. Each thread 1104 includes a respective thread priority. Each thread 1104 in the first data field is sorted with respect to each other thread in the first data field based on a semantic. The semantic is that a thread with a high priority is removed from the run queue before a thread with a lower priority is removed from the run queue. For example, thread 1104-2 having a priority of two (2) is removed for execution before the thread 1104-3 having a priority of four (4) is removed for execution.

The run queue 1100 comprises one or more second data fields 1108 such as field 1108-1 and field 1108-2. A second data field comprises a second plurality of threads 1104. Each thread 1104 in the second data field is sorted with respect to each other thread in the second data field such that a thread with a high priority is removed from the second data field before a thread with a lower priority is removed.

In this implementation an operating system obtains the secondary data field from the sleep queue 300 of Fig. 3. Such a sleep queue and procedures to remove

1 the secondary data field 308 from the sleep queue in a deterministic amount of
2 time is described in greater detail above in reference to data field 308 of Figs.
3 3-10.

4 The second data field 1108 , which corresponds to the data field 308 of Fig.
5 3, includes a root thread such as root thread 1104-2 or root thread 1104-4. A root
6 thread includes a particular priority, and each of the other threads include a
7 respective priority that is a lower priority or an equal priority as compared to the
8 root thread's priority. For example, root thread 1104-2 has a priority of 2, each of
9 the other priorities in the second data field 1108-1 have a respective priority that is
10 less than 2. (In this example, a priority of zero (0) such as node 1104-2 is a
11 highest priority, and a priority of fifteen (15) such as node 1104-5 is a lowest
12 priority).

13 The entire second data field 1108 is associated with the first data field 1102
14 in response to a single operation such as a linked list insert node operation,
15 whereupon the root node such as thread 1104-2 is inserted into the run queue 1100
16 first data field 1102. Linked list insert node operations are well known in the art
17 of computer programming. Because the second data field is associated with the
18 first data field in response to a single operation, the amount of time that it takes to
19 associate any number of threads with the run queue is independent of a number of
20 threads being associated with the run queue.

21 Moreover, because each thread in the data field 1108 is associated with the
22 run queue in response to a single operation, the amount of time that it takes to
23 associate any number of threads with the run queue is determinable in advance of
24 performing the operation. A single operation may include one or more instructions
25 each of which are determinable in advance. Because a processor typically requires

1 a predetermined amount of time to execute any one particular instruction, the
2 amount of time to execute the operation to associate the group of threads with the
3 run queue is deterministic. Because the amount of time is deterministic, it is also
4 often referred to as being "bounded" since one can specify a boundary limit to the
5 amount of time in advance.

6 Only those threads 1104 that are in the first data field's queue 1102 will be
7 removed for execution. A thread of a higher priority will always be removed from
8 the first data queue for execution before a thread of a lower priority is removed
9 from the first data queue for execution. And, threads of equal priority run in a
10 first-in-first-out round-robin fashion. Note that a root thread such a thread 1104-2
11 and/or a thread 1104-4 is part of both a first data queue and a respective second
12 data queue 1108.

13 In this implementation, it is only after a root node has been removed for
14 execution, that a next node coupled to the removed root node is inserted into the
15 first data queue for subsequent removal and execution. For example, only after
16 the root node 1104-2 has been removed for execution, is the next node 1104-6
17 inserted into the first data queue for subsequent removal and execution. The next
18 node is inserted such that it maintains the priority sorting of the first data queue.
19 In this example, the next node has a priority of four (4), thus it would be inserted
20 into the first queue such that it will be removed for execution before node 1104-3
21 is removed for execution. However, had the next node had a lower priority than
22 node 1104-3 (such as a priority of 20), the next node would be inserted into the
23 first queue such that it would not be removed for execution until after the
24 node 1104-5 is removed for execution. More particularly, the next node is inserted
25

1 into the first queue such that it is not removed from the first data queue until after
2 nodes 1104-3 through 1104-5 have been removed for execution.

3 In this manner, after a root node such as thread 1104-2 and/or thread 1104-4
4 is removed from the first data queue 1102, a next node (if there is one) such
5 as node 1104-6 and/or node 1104-9, in effect, becomes a root node of a respective
6 second data queue 1108. This is because the next node effectively becomes the
7 queue's head node.

8 **Exemplary Procedure to Manage a Run Queue**

9
10 Fig. 12 is a flowchart that shows an exemplary procedure 1200 to remove a
11 thread from run queue 1100 of Fig. 11. At block 1210, the procedure removes a
12 thread from the run queue for execution. As discussed above, the thread may be a
13 root node such as root node 1104-2 of Fig. 11 that is attached to one or more other
14 secondary nodes in data field 1108-1, or may be a node that is not attached to any
15 secondary nodes such as root node 1104-1. Accordingly, at block 1212, the
16 procedure determines whether the removed node is attached to a secondary node
17 (e.g., root node 1104-2 is attached to secondary node 1104-6). If not, the procedure
18 ends.

19 Otherwise, if the removed node is attached to a secondary node (block
20 1212), at block 1214, the procedure inserts the secondary node into the run queue
21 in a manner that maintains the priority based semantics of the run queue. For
22 example, if the removed node (block 1210) is node 1104-4 of Fig. 11, the
23 secondary node is node 1104-9 (having a priority of seven (7)). Block 1214 of the
24 procedure then inserts the secondary node 1104-9 into the run queue before node
25 1104-9, which has a priority of fifteen, thereby making the secondary node a root

node, and thereby maintaining the priority based semantics of the run queue. Significantly, block 1214 of the procedure inserts an additional node into the run queue independent of any access to any other queue such as a sleep queue or a wait queue.

Fig. 13 is a flowchart diagram that shows an exemplary procedure 1300 to insert a plurality of threads into a run queue in a determinative amount of time. At block 1310, the procedure associates a plurality of threads with a run queue of Fig. 11 in a determinative amount of time as described above in reference to Fig. 11. At block 1312, the procedure inserts each thread in the associated plurality of threads (block 1310) into the run queue without an additional sleep queue access (see, sleep queue 300 of Fig. 3).

Exemplary Computing Environment

Fig. 14 illustrates an example of a suitable computing environment 1400 wherein an exemplary system and procedure to manage a run queue may be implemented. Exemplary computing environment 1400 is only one example of a suitable computing environment and is not intended to suggest any limitation as to the scope of use or functionality of an exemplary system and procedure to manage a run queue. The computing environment 1400 should not be interpreted as having any dependency or requirement relating to any one or combination of components illustrated in the exemplary computing environment 1400.

The exemplary system and procedure to manage a run queue is operational with numerous other general purpose or special purpose computing system environments or configurations. Examples of well known computing systems, environments, and/or configurations that may be suitable for use with an

exemplary system and procedure to manage a run queue include, but are not limited to, personal computers, server computers, thin clients, thick clients, handheld or laptop devices, multiprocessor systems, microprocessor-based systems, set top boxes, programmable consumer electronics, wireless phones, application specific integrated circuits (ASICs), network PCs, minicomputers, mainframe computers, distributed computing environments that include any of the above systems or devices, and the like.

Exemplary run queue management may be described in the general context of computer-executable instructions, such as program modules, being executed by a computer. Generally, program modules include routines, programs, objects, components, data structures, etc. that perform particular tasks or implement particular abstract data types. Exemplary run queue management may also be practiced in distributed computing environments where tasks are performed by remote processing devices that are linked through a communications network. In a distributed computing environment, program modules may be located in both local and remote computer storage media including memory storage devices.

As shown in Fig. 14, the computing environment 1400 includes a general-purpose computing device in the form of a computer 1410. The components of computer 1410 may include, by are not limited to, one or more processors or processing units 1412, a system memory 1414, and a bus 1416 that couples various system components including the system memory 1414 to the processor 1412.

Bus 1416 represents one or more of any of several types of bus structures, including a memory bus or memory controller, a peripheral bus, an accelerated graphics port, and a processor or local bus using any of a variety of bus

architectures. By way of example, and not limitation, such architectures include Industry Standard Architecture (ISA) bus, Micro Channel Architecture (MCA) bus, Enhanced ISA (EISA) bus, Video Electronics Standards Association (VESA) local bus, and Peripheral Component Interconnects (PCI) bus also known as Mezzanine bus.

Computer 1410 typically includes a variety of computer-readable media. Such media may be any available media that is accessible by computer 1410, and it includes both volatile and non-volatile media, removable and non-removable media.

In Fig. 14, the system memory 1414 includes computer readable media in the form of volatile memory, such as random access memory (RAM) 1420, and/or non-volatile memory, such as read only memory (ROM) 1418. A basic input/output system (BIOS) 1422, containing the basic routines that help to transfer information between elements within computer 1410, such as during start-up, is stored in ROM 1418. RAM 1420 typically contains data and/or program modules that are immediately accessible to and/or presently be operated on by processor 1412.

Computer 1410 may further include other removable/non-removable, volatile/non-volatile computer storage media. By way of example only, Fig. 14 illustrates a hard disk drive 1424 for reading from and writing to a non-removable, non-volatile magnetic media (not shown and typically called a “hard drive”), a magnetic disk drive 1426 for reading from and writing to a removable, non-volatile magnetic disk 1428 (e.g., a “floppy disk”), and an optical disk drive 1430 for reading from or writing to a removable, non-volatile optical disk 1432 such as a CD-ROM, DVD-ROM or other optical media. The hard disk drive 1424,

The drives and their associated computer-readable media provide nonvolatile storage of computer readable instructions, data structures, program modules, and other data for computer 1410. Although the exemplary environment described herein employs a hard disk, a removable magnetic disk 1428 and a removable optical disk 1432, it should be appreciated by those skilled in the art that other types of computer readable media which can store data that is accessible by a computer, such as magnetic cassettes, flash memory cards, digital video disks, random access memories (RAMs), read only memories (ROM), and the like, may also be used in the exemplary operating environment.

A number of program modules may be stored on the hard disk, magnetic disk 1428, optical disk 1432, ROM 1418, or RAM 1420, including, by way of example, and not limitation, an operating system 1438, one or more application programs 1440, other program modules 1442, and program data 1444. Each such operating system 1438, one or more application programs 1440, other program modules 1442, and program data 1444 (or some combination thereof) may include an implementation to manage a run queue.

A user may enter commands and information into computer 1410 through input devices such as keyboard 1446 and pointing device 1448 (such as a “mouse”). Other input devices (not shown) may include a microphone, joystick, game pad, satellite dish, serial port, scanner, or the like. These and other input devices are connected to the processing unit 1412 through a user input interface 1450 that is coupled to bus 1416, but may be connected by other interface and bus structures, such as a parallel port, game port, or a universal serial bus (USB).

1 A monitor 1452 or other type of display device is also connected to bus
2 1416 via an interface, such as a video adapter 1454. In addition to the monitor,
3 personal computers typically include other peripheral output devices (not shown),
4 such as speakers and printers, which may be connected through output peripheral
5 interface 1455.

6 Computer 1410 may operate in a networked environment using logical
7 connections to one or more remote computers, such as a remote computer 1462.
8 Remote computer 1462 may include many or all of the elements and features
9 described herein relative to computer 1410.

10 Logical connections shown in Fig. 14 are a local area network (LAN) 1457
11 and a general wide area network (WAN) 1459. Such networking environments
12 are commonplace in offices, enterprise-wide computer networks, intranets, and the
13 Internet.

14 When used in a LAN networking environment, the computer 1410 is
15 connected to LAN 1457 via network interface or adapter 1466. When used in a
16 WAN networking environment, the computer typically includes a modem 1458 or
17 other means for establishing communications over the WAN 1459. The modem
18 1458, which may be internal or external, may be connected to the system bus 1416
19 via the user input interface 1450 or other appropriate mechanism.

20 Depicted in Fig. 14 is a specific implementation of a WAN via the Internet.
21 Computer 1410 typically includes a modem 1458 or other means for establishing
22 communications over the Internet 1460. Modem 1458, which may be internal or
23 external, is connected to bus 1416 via interface 1450.

24 In a networked environment, program modules depicted relative to the
25 personal computer 1410, or portions thereof, may be stored in a remote memory

1 storage device. By way of example, and not limitation, Fig. 14 illustrates remote
2 application programs 1469 as residing on a memory device of remote computer
3 1462. It will be appreciated that the network connections shown and described are
4 exemplary and other means of establishing a communications link between the
5 computers may be used.

6 **Computer-Executable Instructions**

8 An implementation to manage a run queue may be described in the general
9 context of computer-executable instructions, such as program modules, executed
10 by one or more computers or other devices. Program modules typically include
11 routines, programs, objects, components, data structures, and the like, that perform
12 particular tasks or implement particular abstract data types. The functionality of
13 the program modules typically may be combined or distributed as desired in the
14 various embodiments of Fig. 14.

15 **Computer Readable Media**

17 An implementation to manage a run queue may be stored on or transmitted
18 across some form of computer-readable media. Computer-readable media can be
19 any available media that can be accessed by a computer. By way of example, and
20 not limitation, computer readable media may comprise “computer storage media”
21 and “communications media.”

22 “Computer storage media” include volatile and non-volatile, removable and
23 non-removable media implemented in any method or technology for storage of
24 information such as computer readable instructions, data structures, program
25 modules, or other data. Computer storage media includes, but is not limited to,

1 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,
2 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic
3 tape, magnetic disk storage or other magnetic storage devices, or any other
4 medium which can be used to store the desired information and which can be
5 accessed by a computer.

6 "Communication media" typically embodies computer readable
7 instructions, data structures, program modules, or other data in a modulated data
8 signal, such as carrier wave or other transport mechanism. Communication media
9 also includes any information delivery media.

10 The term "modulated data signal" means a signal that has one or more of its
11 characteristics set or changed in such a manner as to encode information in the
12 signal. By way of example, and not limitation, communication media includes
13 wired media such as a wired network or direct-wired connection, and wireless
14 media such as acoustic, RF, infrared, and other wireless media. Combinations of
15 any of the above are also included within the scope of computer readable media.

16 **Conclusion**

17
18 Although various implementations to manage a sleep queue with
19 deterministic time for system-exclusive access have been described in language
20 specific to structural features and/or methodological operations, it is to be
21 understood that the described subject matter to manage a sleep queue with
22 bounded time for system-exclusive access defined in the appended claims is not
23 necessarily limited to the specific features or operations described. Rather, the
24 specific features and operations are disclosed as preferred forms of implementing
25 the claimed present subject matter.